

1

Statistical thermodynamic foundations

We are interested in a system composed of a biological macromolecule and a number of ligands, of which the solvent can be considered one, under conditions of temperature and pressure of biological relevance. Our goal is to characterize the behavior of the macromolecule in its interaction with the various components of the system and the rules underlying the mutual interference of physical and chemical variables. In this chapter we deal with the statistical thermodynamic foundations of binding processes and the concepts that form the basis of our treatment.

1.1 Postulates and basic ensembles

The physico-chemical properties of a system are defined thermodynamically by a set of macroscopic quantities accessible to experimental measurements (Fermi, 1936; Schrödinger, 1946). If the macromolecule is taken as the system, then all observables reflect properties of the system and the way it is affected by physical and chemical driving forces. A macromolecule in the presence of multiple ligands can be seen as a system existing in a number of distinct energy states, E_1, E_2, \dots, E_r , totally analogous to the energy states of a quantum mechanical system. Each energy state may be coupled to a conformational state of the macromolecule, i.e., a specific arrangement of its secondary, tertiary or quaternary structure. In addition, E_j is specified by the number of ligands bound, or by a particular ligated configuration. Different energy states may group together in the same energy level according to their degeneracy. We consider a prototypic system composed of the macromolecule and the solvent, in contact with a heat bath at constant temperature T , and focus our attention on a particular configuration of the macromolecule (our system) with N water molecules (the solvent) bound to it, as shown in

Figure 1.1. We use the term *bound* in a thermodynamic sense, to distinguish the water molecules in any form of interaction with the system from those belonging to the bulk solvent. In what follows, we make no attempt to discriminate among the various binding modes that characterize protein hydration, whether specific or non-specific. We use the macromolecule in solution to illustrate the salient thermodynamic features of the simplest system of interest. If V is the volume of the macromolecule, then any energy state can be written as $E_j(N, V)$ and the spectrum of energy values E_1, E_2, \dots, E_r denotes the possible alternative states accessible to

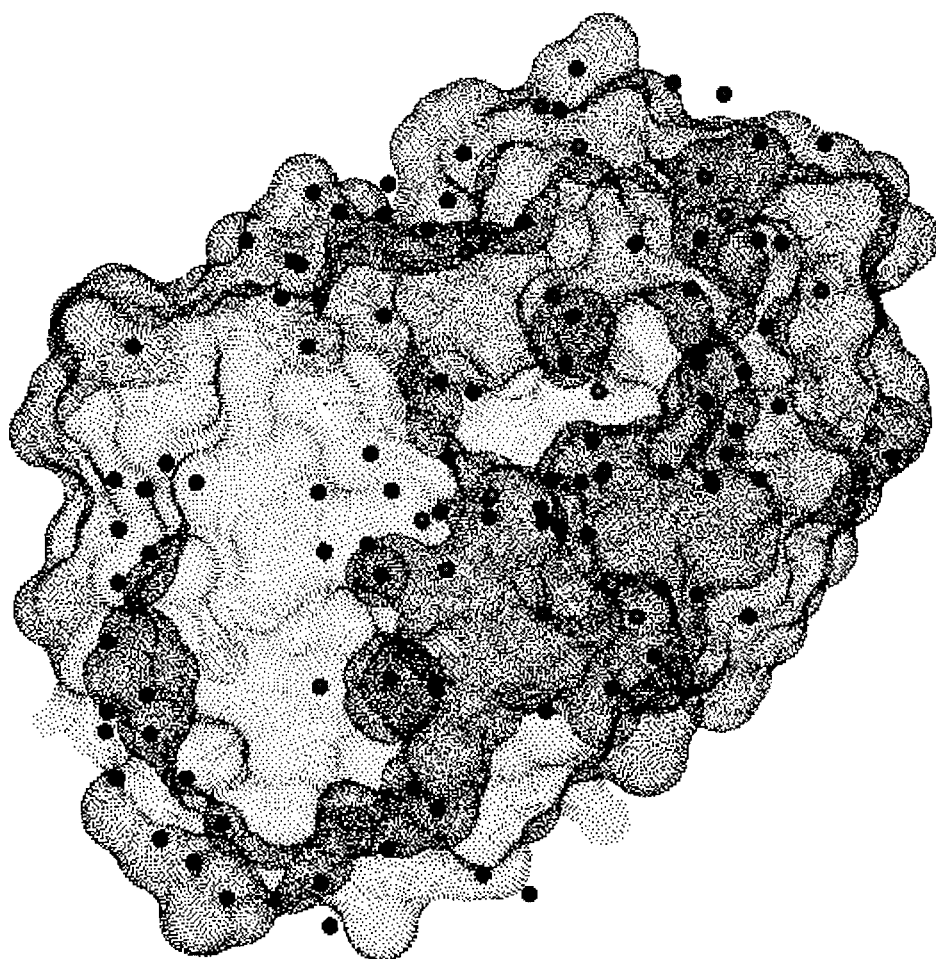


Figure 1.1 Connolly surface of lysozyme, constructed from the crystal structure (Weaver and Matthews, 1987). Water molecules are indicated by solid dots. Some molecules are present as components of the bulk solvent. Other molecules are 'bound' to the accessible surface of the enzyme and make polar contacts with specific residues.

1.1 Postulates and basic ensembles

3

the macromolecule with specified values of V and N . The question that arises at this point concerns the probability of finding the macromolecule in a particular energy state, E_j , when the temperature of the heat bath is held constant. To solve this problem we use the ensemble method introduced by Gibbs (Gibbs, 1902, 1928; Hill, 1960).

Consider an ensemble of a very large number of identical replicas of our system, as depicted in Figure 1.2, with the same value of V and N , and in thermal contact with each other. Each system is in contact with a heat bath at constant T and is allowed to exchange heat, but not water molecules, with the surrounding systems. After equilibrium is reached and T is uniform everywhere, the entire ensemble of systems, or supersystem, is thermally insulated by an adiabatic wall. If Γ is the total number of replicas in the supersystem, then each system can be thought of as being in contact with a heat bath at temperature T formed by the remaining $\Gamma - 1$ systems. The supersystem artificially generated by our replication process is an *isolated* system that cannot exchange heat or matter with the environment, as opposed to each constituent system of the ensemble that is *closed* and isothermal, since it is allowed to exchange heat, but not matter, with the environment. The properties of our original system can be derived from the properties of the supersystem characterized by a

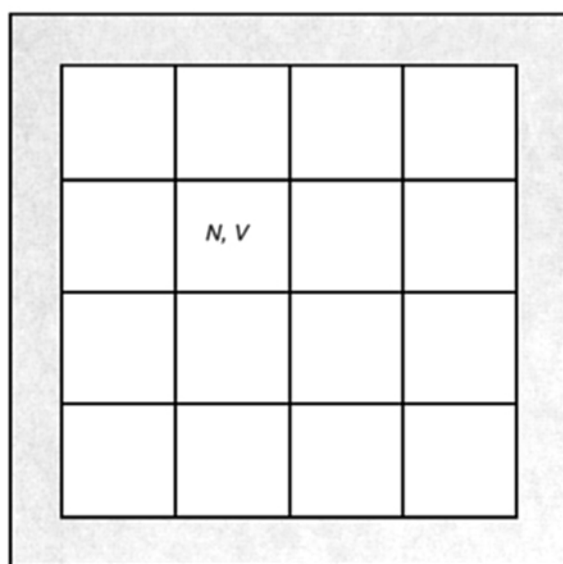


Figure 1.2 Schematic representation of a canonical ensemble. A number of elementary systems with fixed values of V and N are in thermal contact with each other at constant temperature. The supersystem formed by Γ elementary systems is insulated by an adiabatic wall from the environment.

volume ΓV , ΓN water molecules and a total energy E_t . In the thermodynamic limit $\Gamma \rightarrow \infty$, the supersystem is composed of an extremely large number of replicas and the properties of each system separately can be derived as ‘ensemble averages’ taken over the entire collection of replicas. The validity of this assertion is supported by two postulates of Gibbs.

The first postulate states that:

The long time average of a variable in the system at equilibrium is equal to its ensemble average in the thermodynamic limit $\Gamma \rightarrow \infty$, provided the ensemble of systems replicates the system of interest and its environment.

This postulate is intuitively obvious. It simply states the equivalence between averages obtained by monitoring the behavior of the system at equilibrium over a long time scale and those obtained over the ensemble artificially constructed. The time evolution of a variable of interest, say the energy of the system, is plotted in Figure 1.3 as a trajectory whose

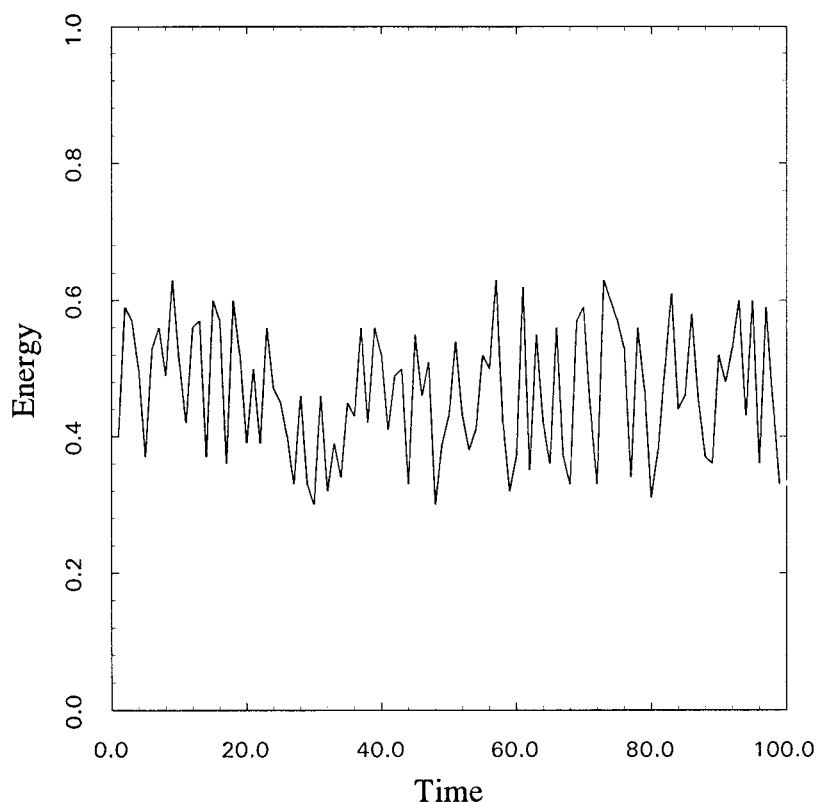


Figure 1.3 Time evolution of the energy (in arbitrary units) of a canonical ensemble at equilibrium. The time average of the energy values gives the thermodynamic measure of the energy E .

time average gives the thermodynamic measure of $E(N, V)$. The contact with the heat bath makes the energy fluctuate and assume any of the allowable values in the spectrum. If we were to monitor the energy for our system at equilibrium, we would obtain a plot similar to that shown in Figure 1.3. The average value of the energy must, of course, be independent of the particular time at which the observation is started. Hence, we could make repeated observations about our system and obtain a bundle of trajectories such as the one shown in Figure 1.3, each of which would yield the same value of the average energy $E(N, V)$. We may think of each trajectory as belonging to a replica of our original system. In the limit where this number becomes arbitrarily large, we can take a snapshot of the bundle of trajectories at any given time and display the energy values for all individual systems, as shown in Figure 1.4. The average of these

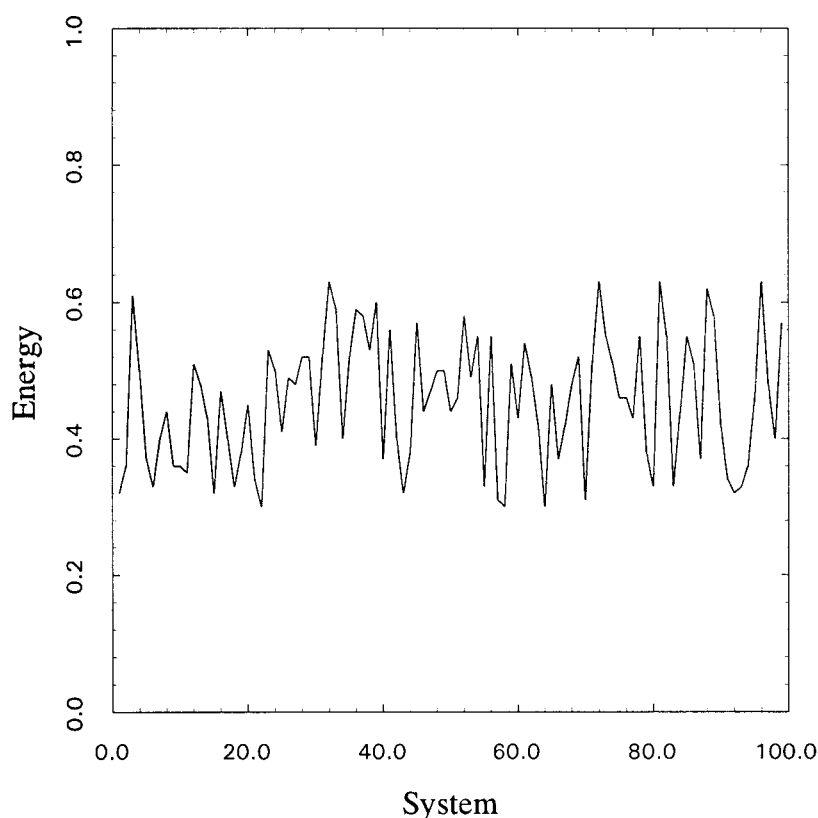


Figure 1.4 Spectrum of energy values for a bundle of trajectories such as the one shown in Figure 1.3, taken at a fixed time, for a canonical ensemble at equilibrium. The abscissa enumerates a representative sample of systems in the supersystem in Figure 1.2. The ensemble average of the energy values is equivalent to the time average derived from Figure 1.3 and gives the thermodynamic measure of the energy E .

values must be independent of the particular time chosen. Sections of the bundle taken at different times must yield energy profiles with the same average value, and this value is evidently the same as that obtained in the long time average of a single trajectory. Hence the first postulate, which establishes a complete equivalence between time averages over the system and ensemble averages independent of time over a supersystem constructed of a large number of replicas of the original system.

In order to determine the properties of our system specifically we need a second postulate, also known as the ‘principle of equal *a priori* probabilities’ for an isolated system, that states:

*In an isolated system for which N , V and the energy E are specified, all energy states have equal *a priori* probabilities.*

All allowable energy states, consistent with the specified values of N , V and E , occur with the same frequency if the behaviour of the system is followed over time or, equivalently, the system spends equal amounts of time in any of the allowable energy states. This important consequence of the two postulates is known as the ‘ergodic hypothesis’ and plays a central conceptual role in statistical mechanics and its mathematical foundations (Tolman, 1938; Khinchin, 1949; Kubo, 1965; de Groot and Mazur, 1984). In an ensemble representative of an isolated system, the probability of occurrence of any energy state is exactly the same for all systems and if any system is chosen at random, it will be in any of the allowable energy states with equal probability. The second postulate as stated above applies only to *isolated* systems, while the first postulate holds quite generally for any equilibrium system.

We are now in the position to find the properties of our original system composed of the macromolecule with a volume V , a number N of water molecules bound, and in contact with a heat bath at constant temperature T . The ensemble of elementary systems shown in Figure 1.2 is, by construction, an isolated system which must obey the Gibbs postulates. The energy of the supersystem is

$$E_t = \sum_{j=1}^r n_j E_j \quad (1.1)$$

where n_j is the number of systems in the energy state E_j . Also,

$$\Gamma = \sum_{j=1}^r n_j \quad (1.2)$$

1.1 Postulates and basic ensembles

7

is an obvious conservation relationship. There is a very large number of distributions, D , of n_j values consistent with eqs (1.1) and (1.2), and specifically there are as many as

$$\Omega(D) = \frac{\Gamma!}{n_1(D)!n_2(D)! \dots n_r(D)!} \quad (1.3)$$

possible ways of realizing the particular distribution, D , of such values. Here $n_j(D)$ denotes the number of systems in state E_j for the particular distributions D . The probability, ψ_j , of finding the macromolecule in a given energy state E_j is the weighted mean of all possible $n_j(D)$ values, divided by the total number of systems, i.e.,

$$\psi_j = \frac{\sum_n n_j(D)\Omega(D)}{\Gamma \sum_n \Omega(D)} \quad (1.4)$$

The distribution of n_j values, $\Omega(D)$, is multinomial and goes into a Gaussian for large Γ by virtue of the De Moivre–Laplace theorem (Wilson, 1911; Feller, 1950). In the limit $\Gamma \rightarrow \infty$, the Gaussian goes into a completely sharp Dirac δ -function centered about the most probable distribution $D = D^*$. Since the δ -function vanishes everywhere except for $D = D^*$, then the simple result $\psi_j = n_j^*/\Gamma$ is obtained from eq (1.4), where $n_j^* = n_j(D^*)$ is the value of n_j in the most probable distribution.

The value of n_j^* can be found by maximizing Ω , or equivalently $\ln \Omega$, subject to the constraints in eqs (1.1) and (1.2). Application of the method of Lagrange's undetermined multipliers (Wilson, 1911) gives

$$\frac{\partial}{\partial n_j} [\ln \Omega(D) - \alpha \Gamma - \beta E_j] = 0 \quad (1.5)$$

where α and β are the undetermined multipliers. Using the Stirling approximation $\ln x! \approx x \ln x - x$ for the factorial terms of $\Omega(D)$ in eq (1.3), and recalling that Γ and E_j are functions of n_j , yields

$$\ln \Gamma - \ln n_j^* - \alpha - \beta E_j = 0 \quad (1.6)$$

Hence,

$$\psi_j = \frac{n_j^*}{\Gamma} = \exp(-\alpha - \beta E_j) \quad (1.7)$$

gives the probability that the macromolecule exists in the energy state E_j defined by V and N . The value of β is $(k_B T)^{-1}$, where k_B is the Boltzmann constant (Tolman, 1938; Hill, 1960). The value of α is

obtained from the conservation relationship (1.2), which embodies the obvious fact that the sum of all ψ_j must equal unity, so that

$$\alpha = \ln \sum_{j=1}^r \exp\left(-\frac{E_j}{k_B T}\right) \quad (1.8)$$

Hence,

$$\psi_j = \frac{\exp\left[-\frac{E_j(N, V)}{k_B T}\right]}{\sum_{j=1}^r \exp\left[-\frac{E_j(N, V)}{k_B T}\right]} \quad (1.9)$$

gives the explicit expression for the probability of the macromolecule existing in the energy state E_j . This quantity decreases exponentially with the value of E_j . The ensemble average over all possible energy states

$$E = \sum_{j=1}^r \psi_j E_j = \langle E \rangle \quad (1.10)$$

defines the energy of the system in the thermodynamic sense.

The closed isothermal system considered in the foregoing analysis is called a *canonical ensemble* and the probability distribution in eq (1.9) is the Boltzmann distribution for a canonical ensemble (Tolman, 1938). The properties of the macromolecule considered as a closed isothermal system, or a canonical ensemble, are completely specified by its volume V , the number of water molecules bound N and the temperature of the heat bath T . Thermodynamically, such a system is an N - V - T ensemble. We now introduce an important quantity that allows us to compute all thermodynamic functions of interest for the system. This is the canonical *partition function* defined as the sum

$$Z(N, V, T) = \sum_{j=1}^r \exp\left[-\frac{E_j(N, V)}{k_B T}\right] \quad (1.11)$$

The partition function enumerates all possible energy states of the system as they appear in the Boltzmann distribution. The relative contribution of each term is weighted exponentially by the Boltzmann factor associated with it. The probability of any energy state, E_j , is given by the ratio between the term containing the value of E_j and the partition function, as seen in eq (1.9). Mathematically, this can be expressed in compact form with the use of the partial derivative

1.1 Postulates and basic ensembles

9

$$\psi_j = -k_B T \frac{\partial \ln Z}{\partial E_j} = \frac{\partial F}{\partial E_j} \quad (1.12)$$

where

$$F(N, V, T) = -k_B T \ln Z(N, V, T) \quad (1.13)$$

is the Helmholtz free energy of the macromolecule and plays the role of the potential associated with the partition function Z . The energy of the macromolecule can be derived in an analogous way as

$$E = k_B T^2 \left(\frac{\partial \ln Z}{\partial T} \right)_{N,V} = -T^2 \left(\frac{\partial \frac{F}{T}}{\partial T} \right)_{N,V} = \left(\frac{\partial \frac{F}{T}}{\partial \frac{1}{T}} \right)_{N,V} \quad (1.14)$$

Since Z is also a function of V and N , it is important to derive the quantities associated with these independent variables. The derivative

$$P = k_B T \left(\frac{\partial \ln Z}{\partial V} \right)_{N,T} = - \left(\frac{\partial F}{\partial V} \right)_{N,T} \quad (1.15)$$

gives the pressure of the system and

$$\mu = -k_B T \left(\frac{\partial \ln Z}{\partial N} \right)_{V,T} = \left(\frac{\partial F}{\partial N} \right)_{V,T} \quad (1.16)$$

is the chemical potential of water.

The change of Z or F with respect to the external conditions subject to experimental control yields information on the quantities associated with them. The energy of the system is obtained as the ‘response function’ to a change in temperature. Likewise, the pressure and chemical potential of water are obtained as responses to changes of V and N respectively. From the definition of F in eq (1.13) it also follows that

$$\left(\frac{\partial F}{\partial T} \right)_{N,V} = -k_B \ln Z - k_B T \left(\frac{\partial \ln Z}{\partial T} \right)_{N,V} = \frac{F - E}{T} = -S \quad (1.17)$$

The entropy, S , of the macromolecule is derived from the change of the Helmholtz free energy as a function of temperature. An explicit expression for S is obtained from the definition of ψ_j as follows. Taking the logarithm of ψ_j one has

$$\ln \psi_j = -\frac{E_j}{k_B T} - \ln Z \quad (1.18)$$

Multiplying both members by ψ_j and summing over all values of j leads to

$$\frac{F - E}{k_B T} = \sum_{j=1}^r \psi_j \ln \psi_j \quad (1.19)$$

Hence,

$$S = -k_B \sum_{j=1}^r \psi_j \ln \psi_j = -\frac{F - E}{T} \quad (1.20)$$

provides an important definition of the entropy in terms of the probability of occurrence of the energy states, and the thermodynamic potentials E and F . Eqs (1.15)–(1.17) also allow for a definition of the thermodynamic potential F in differential, or Pfaffian form (de Heer, 1986) as follows

$$\begin{aligned} dF &= \left(\frac{\partial F}{\partial N} \right)_{V,T} dN + \left(\frac{\partial F}{\partial V} \right)_{N,T} dV + \left(\frac{\partial F}{\partial T} \right)_{N,V} dT \\ &= \mu dN - P dV - S dT \end{aligned} \quad (1.21)$$

This Pfaffian form summarizes the properties of the macromolecule as a canonical ensemble. The integral form of F is

$$F = E - TS \quad (1.22)$$

and follows directly from eq (1.20).

Having defined the properties of our original closed isothermal system, we turn to the supersystem itself which is an isolated system for which N , V and the total energy E are fixed. Consider all systems in the supersystem that belong to the same energy state E_j , group and surround them by an adiabatic wall. This yields an ensemble of systems with the same value of N , V and $E = E_j$. Such an ensemble is called *microcanonical* and differs from the canonical ensemble insofar as all systems in the ensemble have the same energy. Application of the same arguments developed for the canonical ensemble leads to the conclusion that ψ_j must be independent of j in the microcanonical ensemble. This is because the microcanonical ensemble is a degenerate canonical ensemble where all systems have the same energy $E = E_1 = E_2 = \dots = E_r$. We also know, by construction, that the microcanonical ensemble is a subensemble of the canonical ensemble in Figure 1.2, where all systems with the same energy have been grouped together and isolated. Hence, the degeneracy $\Omega(D^*)$ of the most probable distribution D^* in a microcanonical ensemble with energy $E = E_j$ coincides with the value of n_j^* in the canonical ensemble. The value of the probability ψ for the energy E is simply $1/\Omega(D^*)$ and is the same for all systems of the ensemble since they belong to the same energy level and there are a total of $\Omega(D^*)$ such systems. This follows directly from the Boltzmann distribution in eq (1.9) by letting $E_j = E$ for all $j = 1, 2, \dots, \Omega(D^*)$. The partition function of the macromolecule as a microcanonical ensemble is given by Ω and is completely defined once N , V